

PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2003-157271

(43)Date of publication of application : 30.05.2003

(51)Int.Cl.

G06F 17/30

(21)Application number : 2001-355150

(71)Applicant : MITSUBISHI ELECTRIC CORP

(22)Date of filing : 20.11.2001

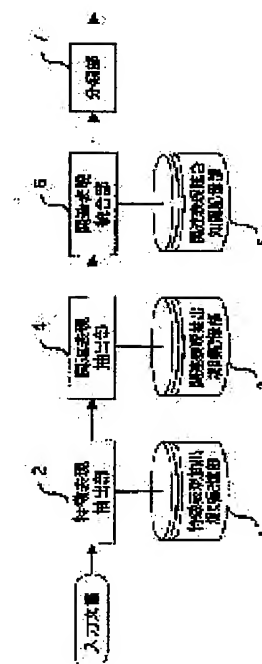
(72)Inventor : NAGAI AKITO
TAKAYAMA YASUHIRO
SUZUKI KATSUSHI

(54) DEVICE AND METHOD FOR MINING TEXT

(57)Abstract:

PROBLEM TO BE SOLVED: To obtain a text mining device and a method having such a general usability that can analyze even Web pages as analysis objects.

SOLUTION: The text mining device is provided with a feature expression extraction rule storage part 1, a feature expression extraction part 2 for extracting a 1st description unit including intended extraction expression as a feature expression extraction result and outputting the extraction result together with the extraction intention, a related expression extraction rule storage part 3 for allowing a relation name of each 1st description unit to correspond to relation expression, and a related expression extraction part 4 for extracting a 2nd description unit from the feature expression extraction result as a related expression extraction result and outputting the extracted result together with the relation name. The device is provided also with a related expression integrated knowledge storage part 5 for integrating the feature expression and the related expression on the basis of relation of each 1st description unit, a related expression integration part 6 for outputting a related expression integration result and a classification part 7 for classifying the related expression integration result in accordance with the contents of the extraction intention.



(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開2003-157271

(P2003-157271A)

(43)公開日 平成15年5月30日(2003.5.30)

(51)Int.Cl. ⁷	識別記号	F I	テーマコード*(参考)
G 0 6 F 17/30	2 2 0	G 0 6 F 17/30	2 2 0 Z 5 B 0 7 5
	1 7 0		1 7 0 A
	3 3 0		3 3 0 C

審査請求 未請求 請求項の数16 O L (全 12 頁)

(21)出願番号 特願2001-355150(P2001-355150)

(22)出願日 平成13年11月20日(2001.11.20)

(71)出願人 000006013

三菱電機株式会社

東京都千代田区丸の内二丁目2番3号

(72)発明者 永井 明人

東京都千代田区丸の内二丁目2番3号 三

菱電機株式会社内

(72)発明者 高山 泰博

東京都千代田区丸の内二丁目2番3号 三

菱電機株式会社内

(74)代理人 100057874

弁理士 曾我 道照 (外6名)

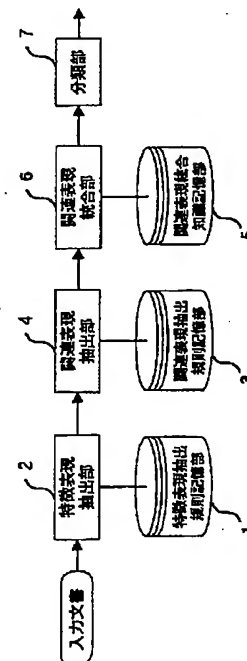
最終頁に続く

(54)【発明の名称】 テキストマイニング装置および方法

(57)【要約】

【課題】 Webページも含めて分析対象にできるような汎用性を有するテキストマイニング装置および方法を得る。

【解決手段】 特徴表現抽出規則記憶部1と、意図抽出表現を含む第1の記述単位を抽出して特徴表現抽出結果とし、抽出意図とともに出力する特徴表現抽出部2と、第1の記述単位毎の関係名と関係表現とを対応付ける関連表現抽出規則記憶部3と、特徴表現抽出結果から第2の記述単位を抽出して関連表現抽出結果とし、関係名とともに出力する関連表現抽出部4と、特徴表現と関連表現とを第1の記述単位毎の関係に基づいて統合する関連表現統合知識記憶部5と、関連表現統合結果を出力する関連表現統合部6と、関連表現統合結果を抽出意図の内容に応じて分類する分類部7とを備えた。



【特許請求の範囲】

【請求項1】 入力文書に関して、抽出意図と意図抽出表現とを対応付けるための特徴表現抽出規則を記憶する特徴表現抽出規則記憶部と、
前記入力文書から、前記意図抽出表現を含む第1の記述単位を抽出して特徴表現とし、前記特徴表現を含む解析結果を特徴表現抽出結果として、前記抽出意図とともに出力する特徴表現抽出部と、
前記第1の記述単位毎の関係に関し、関係名と前記関係を表わす関係表現とを対応付けるための関連表現抽出規則を記憶する関連表現抽出規則記憶部と、
前記特徴表現抽出結果から、前記関係表現を含む第2の記述単位を抽出して関連表現とし、前記関連表現を含む解析結果を関連表現抽出結果として、前記関係名とともに出力する関連表現抽出部と、
前記特徴表現と前記関連表現とを前記第1の記述単位毎の関係に基づいて統合するための関連表現統合知識を記憶する関連表現統合知識記憶部と、
前記関連表現統合知識を参照して前記特徴表現と前記関連表現とを統合し、関連表現統合結果として出力する関連表現統合部と、
関連表現統合結果を前記抽出意図の内容に応じて分類する分類部と
を備えたテキストマイニング装置。

【請求項2】 前記特徴表現抽出部は、前記入力文書を分割して文集合とし、前記意図抽出表現の抽出対象となる第1の記述単位を文とすることを特徴とする請求項1に記載のテキストマイニング装置。

【請求項3】 前記特徴表現抽出部は、前記入力文書における各記述の位置情報を参照し、一定範囲の位置に存在する記述を抽出対象とすることを特徴とする請求項1または請求項2に記載のテキストマイニング装置。

【請求項4】 前記特徴表現抽出部は、文書の構造情報により前記入力文書を分割してサブ文書集合とし、前記サブ文書集合内の各サブ文書毎に前記特徴表現抽出結果を得ることを特徴とする請求項1から請求項3までのいずれかに記載のテキストマイニング装置。

【請求項5】 前記特徴表現抽出部は、文書の構造情報を参照して、前記サブ文書集合の中から、抽出対象とするサブ文書を選択することを特徴とする請求項4に記載のテキストマイニング装置。

【請求項6】 分析の目的とする対象名に関する情報を記憶する対象名辞書と、
前記対象名を含む第3の記述単位を抽出し、前記対象名とともに出力する対象名抽出部とを備え、
前記関連表現統合部は、統合対象となる第1の記述単位である前記特徴表現抽出結果と前記関連表現抽出結果との中から、前記対象名が抽出された第3の記述単位に前記対象名を付与して出力し、
前記分類部は、前記対象名毎に分類して提示することを

特徴とする請求項1から請求項5までのいずれかに記載のテキストマイニング装置。

【請求項7】 前記関連表現統合部は、前記特徴表現抽出結果と前記関連表現抽出結果との統合判定に際して、前記対象名を利用することを特徴とする請求項6に記載のテキストマイニング装置。

【請求項8】 前記分類部は、前記関連表現統合部により得られた前記関連表現統合結果を事例として、前記事例のクラスタリングを行い、問題解決木を構成して出力することを特徴とする請求項7に記載のテキストマイニング装置。

【請求項9】 入力文書に関して、抽出意図と意図抽出表現とを対応付けるための特徴表現抽出規則を記憶する特徴表現抽出規則記憶ステップと、
前記入力文書から、前記意図抽出表現を含む第1の記述単位を抽出して特徴表現とし、前記特徴表現を含む解析結果を特徴表現抽出結果として、前記抽出意図とともに出力する特徴表現抽出ステップと、前記第1の記述単位毎の関係に関し、関係名と前記関係を表わす関係表現とを対応付けるための関連表現抽出規則を記憶する関連表現抽出規則記憶ステップと、
前記特徴表現抽出結果から、前記関係表現を含む第2の記述単位を抽出して関連表現とし、前記関連表現を含む解析結果を関連表現抽出結果として、前記関係名とともに出力する関連表現抽出ステップと、
前記特徴表現と前記関連表現とを前記第1の記述単位毎の関係に基づいて統合するための関連表現統合知識を記憶する関連表現統合知識記憶ステップと、
前記関連表現統合知識を参照して前記特徴表現と前記関連表現とを統合し、関連表現統合結果として出力する関連表現統合ステップと、
前記関連表現統合結果を前記抽出意図の内容に応じて分類する分類ステップとを備えたテキストマイニング方法。

【請求項10】 前記特徴表現抽出ステップは、前記入力文書を分割して文集合とし、前記意図抽出表現の抽出対象となる第1の記述単位を文とすることを特徴とする請求項9に記載のテキストマイニング方法。

【請求項11】 前記特徴表現抽出ステップは、前記入力文書における各記述の位置情報を参照し、一定範囲の位置に存在する記述を抽出対象とすることを特徴とする請求項9または請求項10に記載のテキストマイニング方法。

【請求項12】 前記特徴表現抽出ステップは、文書の構造情報により前記入力文書を分割してサブ文書集合とし、前記サブ文書集合内の各サブ文書毎に前記特徴表現抽出結果を得ることを特徴とする請求項9から請求項11までのいずれかに記載のテキストマイニング方法。

【請求項13】 前記特徴表現抽出ステップは、文書の構造情報を参照して、前記サブ文書集合の中から、抽出

対象とするサブ文書を選択することを特徴とする請求項12に記載のテキストマイニング方法。

【請求項14】 分析の目的とする対象名に関する情報を記憶する対象名辞書から、前記対象名を含む第3の記述単位を抽出し、前記対象名とともに出力する対象名抽出ステップを備え、

前記関連表現統合ステップは、統合対象となる第1の記述単位である前記特徴表現抽出結果と前記関連表現抽出結果との中から、前記対象名が抽出された第3の記述単位に前記対象名を付与して出力し、

前記分類ステップは、前記対象名毎に分類して提示することを特徴とする請求項9から請求項13までのいずれかに記載のテキストマイニング方法。

【請求項15】 前記関連表現統合ステップは、前記特徴表現抽出結果と前記関連表現抽出結果との統合判定に際して、前記対象名を利用することを特徴とする請求項14に記載のテキストマイニング方法。

【請求項16】 前記分類ステップは、前記関連表現統合ステップにより得られた前記関連表現統合結果を事例として、前記事例のクラスタリングを行い、問題解決木を構成して出力することを特徴とする請求項15に記載のテキストマイニング方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】この発明は、入力される文書中のテキストを解析して特徴的な表現を抽出するテキストマイニング装置および方法に関し、特に、抽出された表現と関係を持つ関連表現をさらに抽出して関係付けることのできるテキストマイニング装置および方法に関するものである。

【0002】

【従来の技術】近年、インターネット上のWebページや電子メール、または蓄積されたテキストデータなどの電子化文書は、非常に勢いで増加しており、人手で全文書の内容を読んで調査することは、最早、不可能になってきている。

【0003】そこで、これらの大量の文書から必要な情報を検索して、内容に応じて分類したり、業務上の分析に役立つ情報を抽出するという技術が、ますます要求されつつある。

【0004】従来より、上記要求に応える装置として、たとえば特開2001-101199号公報（以下、「文献1」という）に記載された「文書処理装置」、または、特開2001-147937号公報（以下、「文献2」という）に記載された「業務支援システム」などが提案されている。

【0005】文献1においては、製品のサポート窓口に送られてきたユーザの電子メールを分析の対象としており、製品の不具合を指摘する電子メールの内容を形態素解析し、形態素を概念に結び付けるための概念分類体系

（シソーラス）を用いて、形態素間の係り受け関係に基づいて、概念間の関係を抽出している。

【0006】また、文献2においては、営業報告書を分析の対象としており、製品の売れ行き情報に関する概念と表現を定義した概念定義辞書を用いて、営業報告内容から概念を抽出し、接続表現と概念関係との対応を定義した接続関係定義辞書を用いて概念間の関係を抽出している。

【0007】また、文献2に関連した特開2001-184351号公報（以下、「文献3」という）に記載された「文書情報抽出装置および文書分類装置」においては、文献2で用いられる概念定義辞書のメンテナンスを容易にするために、辞書登録すべき候補の一覧を提示するための技術が開示されている。

【0008】また、特開2000-357170号公報（以下、「文献4」という）に記載された「文書の参照理由を用いて情報検索を行う装置」においては、技術論文のように、文献間に「参照/被参照関係」が存在する場合に、参照理由を同定して、参照理由を用いて文書を検索する装置が開示されている。

【0009】さらに、立石健二、石黒義英および福島俊一による情報処理学会研究会資料の「インターネットからの評判情報検索」（自然言語処理144-11、第75～82頁、2001年7月16日）（以下、「文献5」という）においては、インターネットのWebページから商品の評価表現を抽出する技術が開示されている。

【0010】図10は上記文献5に記載されたシステム構成を示すブロック構成図である。図10において、10は商品カテゴリに対応した単語を指定する評価表現辞書、11はクローラ12および索引ファイル13を有するWebページ収集部である。

【0011】14は商品名を入力情報とする意見収集部であり、評価表現辞書10およびWebページ収集部11に関連した近接演算処理部15と、近接演算処理部15の出力側に挿入された適正值判定処理部16とを有する。

【0012】17は意見収集部14の出力情報を分析する分類・分析部であり、肯定・否定分類処理部18を有する。分類・分析部17による処理情報は、検索結果として出力される。

【0013】図10に示した従来システムにおいて、評価表現辞書10およびWebページ収集部11に関連した意見収集部14と分類・分析部17とによる評価表現の抽出動作は、以下のようになる。

【0014】まず、意見収集部14内の近接演算処理部15は、Webページ収集部11に商品名を入力し、索引ファイル13により得られたインターネット全文検索結果から、別途指定された商品カテゴリに対応する評価表現辞書10を用いて、評価表現の単語を含む文章を抽

出する。

【0015】たとえば、商品名が「モバイルギア」であって、商品カテゴリが「コンピュータ」であった場合、評価表現辞書10に記憶された「速い」、「重い」、「壊れやす」などの単語や、「好き」、「いい」、「勧め」などの単語を含む文章が、全文検索結果から抽出される。

【0016】次に、分類・分析部17内の肯定・否定分類処理部18は、意見収集部14により抽出された文章に対して、肯定／否定の分類を行い、商品名に対する評価表現の検索結果として出力する。

【0017】

【発明が解決しようとする課題】従来のテキストマイニング装置および方法は以上のように、文献1および文献2に記載された装置および方法によれば、業務内容に関する概念およびシソーラスを業務内容や製品毎に構築する必要があることから、シソーラスの構築に非常に手間がかかるので、適用する業務内容の変更を容易に行うことができないという問題点があった。

【0018】また、上記文献1および文献2の装置および方法では、広範囲な概念を含み、且つ新規の概念や新語も発生するWebページのような文書内容に対しては、あらかじめ概念を定義することができないので、適用することができないという問題点があった。

【0019】また、文献2に記載された装置および方法によれば、概念間の関係を抽出する処理において、関係を表わす表現（「そこで」、「しかし」など）を文書中より抽出し、これらの表現に対応付けられた関係IDを同定した後に、関係IDの前後に最も近接する概念を関係付けており、概念の関係を隣接する二項関係として扱っているため、抽出すべき概念間の関係を取りこぼしてしまうという問題点があった。

【0020】たとえば、上記文献2の装置および方法では、「商品Aは売れている。サンプルを配布したため。しかし、値段が高いという不満がある。」という文書の場合、「値段が高い」という概念は、商品Aに関する「売れ行きが良い」という概念に対して逆接に関係付けられるべきであるが、最も隣接する「サンプルを配布したため。」に関係付けられてしまうという問題点があった。

【0021】また、文献3に記載された装置および方法によれば、辞書登録の候補を提示することにより辞書のメンテナンスの容易化を図っているものの、登録の候補は、概念に対応付けられる表現であって、概念自体が固定の定義となっているので、あらかじめ定義した概念の範囲内における辞書メンテナンスのみに限定されるという問題点があった。

【0022】また、文献4に記載された装置および方法によれば、参照理由の同定に用いてされる「参照特徴－参照理由対応表」において、参照特徴として本文の記述

表現そのものを用いているので、参照理由を表わす表現を全て定義する必要があり、語彙、構文、時制などのバリエーションを考慮すると、参照特徴（定義すべき表現）の数が組合わせ爆発を生じてしまい現実的ではない。うえ、参照特徴を抽出する手段（言語解析部）に関して具体的な実現方法の記載がなく、実現性に乏しいという問題点があった。

【0023】さらに、文献5に記載された装置および方法によれば、抽出対象が評価表現（意見）のみであり、その根拠や理由および原因などの有益な情報が文書中に存在しても抽出することができないという問題点があった。

【0024】この発明は上記のような問題点を解決するためになされたもので、概念やシソーラスを用いずに表層的な表現に着目して表現抽出を行うことにより、Webページも含めて分析対象にできるような汎用性を有するテキストマイニング装置および方法を得ることを目的とする。

【0025】また、この発明は、抽出対象となる表現A（意見、クレーム、障害報告など）に加えて、抽出した表現Aと関係（根拠や理由、原因、例示など）を有する表現Bを抽出して、表現A、Bを関連付けて提示することにより、有益な情報の発見や、それらの関係の分析を支援することのできるテキストマイニング装置および方法を得ることを目的とする。

【0026】

【課題を解決するための手段】この発明に係るテキストマイニング装置は、入力文書に関して、抽出意図と意図抽出表現とを対応付けるための特徴表現抽出規則を記憶する特徴表現抽出規則記憶部と、入力文書から、意図抽出表現を含む第1の記述単位を抽出して特徴表現とし、特徴表現を含む解析結果を特徴表現抽出結果として、抽出意図とともに出力する特徴表現抽出部と、第1の記述単位毎の関係に関し、関係名と関係を表わす関係表現とを対応付けるための関連表現抽出規則を記憶する関連表現抽出規則記憶部と、特徴表現抽出結果から、関係表現を含む第2の記述単位を抽出して関連表現とし、関連表現を含む解析結果を関連表現抽出結果として、関係名とともに出力する関連表現抽出部と、特徴表現と関連表現とを第1の記述単位毎の関係に基づいて統合するための関連表現統合知識を記憶する関連表現統合知識記憶部と、関連表現統合知識を参照して特徴表現と関連表現とを統合し、関連表現統合結果として出力する関連表現統合部と、関連表現統合結果を抽出意図の内容に応じて分類する分類部とを備えたものである。

【0027】また、この発明に係るテキストマイニング装置の特徴表現抽出部は、入力文書を分割して文集合とし、意図抽出表現の抽出対象となる第1の記述単位を文とするものである。

【0028】また、この発明に係るテキストマイニング

装置の特徴表現抽出部は、入力文書における各記述の位置情報を参照し、一定範囲の位置に存在する記述を抽出対象とするものである。

【0029】また、この発明に係るテキストマイニング装置の特徴表現抽出部は、文書の構造情報により入力文書を分割してサブ文書集合とし、サブ文書集合内の各サブ文書毎に特徴表現抽出結果を得るものである。

【0030】また、この発明に係るテキストマイニング装置の特徴表現抽出部は、文書の構造情報を参照して、サブ文書集合の中から、抽出対象とするサブ文書を選択するものである。

【0031】また、この発明に係るテキストマイニング装置は、分析の目的とする対象名に関する情報を記憶する対象名辞書と、対象名を含む第3の記述単位を抽出し、対象名とともに出力する対象名抽出部とを備え、関連表現統合部は、統合対象となる第1の記述単位である特徴表現抽出結果と関連表現抽出結果との中から、対象名が抽出された第3の記述単位に対象名を付与して出力し、分類部は、対象名毎に分類して提示するものである。

【0032】また、この発明に係るテキストマイニング装置の関連表現統合部は、特徴表現抽出結果と関連表現抽出結果との統合判定に際して、対象名を利用するものである。

【0033】また、この発明に係るテキストマイニング装置の分類部は、関連表現統合部により得られた関連表現統合結果を事例として、事例のクラスタリングを行い、問題解決木を構成して出力するものである。

【0034】また、この発明に係るテキストマイニング方法は、入力文書に関して、抽出意図と意図抽出表現とを対応付けるための特徴表現抽出規則を記憶する特徴表現抽出規則記憶ステップと、入力文書から、意図抽出表現を含む第1の記述単位を抽出して特徴表現とし、特徴表現を含む解析結果を特徴表現抽出結果として、抽出意図とともに出力する特徴表現抽出ステップと、第1の記述単位毎の関係に関し、関係名と関係を表わす関係表現とを対応付けるための関連表現抽出規則を記憶する関連表現抽出規則記憶ステップと、特徴表現抽出結果から、関係表現を含む第2の記述単位を抽出して関連表現とし、関連表現を含む解析結果を関連表現抽出結果として、関係名とともに出力する関連表現抽出ステップと、特徴表現と関連表現とを第1の記述単位毎の関係に基づいて統合するための関連表現統合知識を記憶する関連表現統合知識記憶ステップと、関連表現統合知識を参照して特徴表現と関連表現とを統合し、関連表現統合結果として出力する関連表現統合ステップと、関連表現統合結果を抽出意図の内容に応じて分類する分類ステップとを備えたものである。

【0035】また、この発明に係るテキストマイニング方法の特徴表現抽出ステップは、入力文書を分割して文

集合とし、意図抽出表現の抽出対象となる第1の記述単位を文とするものである。

【0036】また、この発明に係るテキストマイニング方法の特徴表現抽出ステップは、入力文書における各記述の位置情報を参照し、一定範囲の位置に存在する記述を抽出対象とするものである。

【0037】また、この発明に係るテキストマイニング方法の特徴表現抽出ステップは、文書の構造情報により入力文書を分割してサブ文書集合とし、サブ文書集合内の各サブ文書毎に特徴表現抽出結果を得るものである。

【0038】また、この発明に係るテキストマイニング方法の特徴表現抽出ステップは、文書の構造情報を参照して、サブ文書集合の中から、抽出対象とするサブ文書を選択するものである。

【0039】また、この発明に係るテキストマイニング方法は、分析の目的とする対象名に関する情報を記憶する対象名辞書から、対象名を含む第3の記述単位を抽出し、対象名とともに出力する対象名抽出ステップを備え、関連表現統合ステップは、統合対象となる第1の記述単位である特徴表現抽出結果と関連表現抽出結果との中から、対象名が抽出された第3の記述単位に対象名を付与して出力し、分類ステップは、対象名毎に分類して提示するものである。

【0040】また、この発明に係るテキストマイニング方法の関連表現統合ステップは、特徴表現抽出結果と関連表現抽出結果との統合判定に際して、対象名を利用するものである。

【0041】また、この発明に係るテキストマイニング方法の分類ステップは、関連表現統合ステップにより得られた関連表現統合結果を事例として、事例のクラスタリングを行い、問題解決木を構成して出力するものである。

【0042】

【発明の実施の形態】実施の形態1. 以下、図面を参照しながら、この発明の実施の形態1について詳細に説明する。図1はこの発明の実施の形態1によるテキストマイニング装置を示すブロック構成図である。

【0043】図1において、1は特徴表現抽出規則記憶部であり、入力文書に関して、抽出意図と意図抽出表現とを対応付けるための特徴表現抽出規則を記憶する。2は特徴表現抽出規則記憶部1と関連した特徴表現抽出部であり、取り込まれた入力文書から、意図抽出表現を含む第1の記述単位を抽出して特徴表現とし、特徴表現を含む解析結果を特徴表現抽出結果として抽出意図とともに出力する。

【0044】3は関連表現抽出規則記憶部であり、第1の記述単位毎の関係に関し、関係名と関係を表わす関係表現とを対応付けるための関連表現抽出規則を記憶する。4は関連表現抽出規則記憶部3と関連した関連表現抽出部であり、特徴表現抽出結果から関係表現を含む第

2の記述単位を抽出して関連表現とし、関連表現を含む解析結果を関連表現抽出結果として関係名とともに出力する。

【0045】5は関連表現統合知識記憶部であり、関連表現抽出部4から取得される特徴表現と関連表現とを、第1の記述単位毎の関係に基づいて統合するための関連表現統合知識を記憶する。

【0046】6は関連表現統合知識記憶部5と関連した関連表現統合部であり、関連表現統合知識を参照して特徴表現と関連表現とを統合し、関連表現統合結果として出力する。7は関連表現統合結果を内容に応じて分類する分類部である。

【0047】次に、図2のフローチャートおよび図3～図8の説明図を参照しながら、図1に示したこの発明の実施の形態1により実行されるテキストマイニング処理動作について説明する。

【0048】図2において、まず、特徴表現抽出部2に入力文書が取り込まれる(ステップS201)。なお、入力文書は、たとえば、インターネット上に存在するWebページや、電子メール、記憶装置や記憶媒体に記録された文書ファイルのデータなどであり、電子化されたテキストであれば何でもよい。

【0049】続いて、特徴表現抽出部2は、入力文書から記述単位Uが取得可能か否かを判定し(ステップS202)、記述単位Uが取得可能(すなわち、YES)と判定されれば、後述の判定ステップS203に進み、記述単位Uが取得不可能(すなわち、NO)と判定されれば、後述のステップS206に進む。

【0050】なお、記述単位Uは、入力文書中のテキストの一部(または、全て)であってもよく、特徴表現抽出部2は、たとえば、句点、疑問符、改行などの字句情報から一文単位を切り出して記述単位Uとするか、または、入力文書中の単語の位置情報を参照し、一定範囲の位置に存在する記述を切り出して記述単位Uとする。

【0051】または、特徴表現抽出部2は、HTML文書のような構造を持った文書に対してタグ情報を参照し、たとえばインターネットの掲示板であれば、個別の記事毎に入力文書を分割してサブ文書とする。

【0052】さらに、特徴表現抽出部2は、タグ情報を参照して、テキスト情報が含まれている記述を選択して記述単位Uとする。こうして切り出された記述単位Uは、一時的なバッファTに格納される。

【0053】ステップS202において、記述単位Uが取得可能(すなわち、YES)と判定された場合、特徴表現抽出部2は、取得した記述単位Uに対して、特徴表現抽出規則記憶部1に記憶された抽出意図と意図抽出表現とを参照し、記述単位Uに意図抽出表現が含まれているか否かを判定する(ステップS203)。

【0054】記述単位Uに意図抽出表現が含まれている(すなわち、YES)と判定されれば、後述のステップ

S204に進み、記述単位Uに意図抽出表現が含まれていない(すなわち、NO)と判定されれば、ステップS202に戻る。

【0055】たとえば、抽出意図が「クレーム」、「要望」、「賞賛」の場合、特徴表現抽出規則記憶部1内の規則は、図3のように定義される。図3において、特徴表現抽出規則は、各抽出意図毎の意図抽出表現として対応付けられ、記述単位U内の単語の共起パターンが、「1.0」、「0.8」などの重み付きで表現された形式で表される。

【0056】ステップS203において、意図抽出表現が記述単位Uに含まれる(すなわち、YES)と判定された場合、この記述単位Uを特徴表現として抽出し、抽出された記述単位Uに抽出意図のラベルを付与する(ステップS204)。また、記述単位Uを特徴表現リストに追加して(ステップS205)、ステップS202に戻る。

【0057】特徴表現リストは、たとえば、図4に示すような情報を持つリストである。図4において、「A社のサポート受付は対応が悪い」など特徴表現に対して、「クレーム」などの抽出意図、および、「21-36(文字)」などの「開始-終了」の位置が対応付けられ、さらに、「1.0」などの重みが付けられる。

【0058】上記特徴表現リストの作成処理(ステップS202～S205)は、記述単位Uが取得できる限り、繰り返し実行される。その後、取得すべき記述単位Uが存在しなくなった(ステップS202でNOと判定された)時点で、特徴表現抽出部2は、特徴表現リストの作成処理を終了し、特徴表現リスト(図4参照)と、一時的なバッファTに格納されている記述単位群(たとえば、図5参照)とを、関連表現抽出部4に出力する。

【0059】これにより、関連表現抽出部4は、上記処理ステップS202～S205で得られた特徴表現リストを読み込み(ステップS206)、各特徴表現に関して、以下の処理を実行する。

【0060】まず、特徴表現リストが空か否かを判定し(ステップS207)、特徴表現リストが空である(すなわち、YES)と判定されれば、後述のステップS212に進む。

【0061】また、特徴表現リストが空でない(すなわち、NO)と判定されれば、関連表現抽出部4は、特徴表現の位置情報を参照し、一時的なバッファTに格納されている記述単位群の中から、特徴表現から一定範囲に存在する記述単位Rを取得する(ステップS208)。

【0062】たとえば、図5に示すような記述単位群の場合、番号「1」の記述単位が抽出意図「賞賛」の特徴表現として抽出されているものとする。また、一定範囲としては、特徴表現の終了位置から後方「80文字」の範囲が設定されているものとする。

【0063】図5に示した記述単位群の場合、番号

「1」の記述単位に関する位置情報（開始～終了）を参照すると「1-26」であり、一定範囲「80文字」内に存在する記述単位Rは、番号「2」～「4」となる。

【0064】次に、取得された記述単位Rに対して、関連表現抽出部4は、関連情報抽出規則記憶部3に記憶された「関係名」と「関係を表わす関係表現」とを参照し、記述単位Rに関係表現が含まれているか否かを判定する（ステップS209）。

【0065】ステップS209において、記述単位Rに関係表現が含まれている（すなわち、YES）と判定されれば、後述のステップS210に進み、記述単位Rに関係表現が含まれていない（すなわち、NO）と判定されれば、ステップS207に戻る。

【0066】関連情報抽出規則記憶部3内の規則は、たとえば、図6のように定義されている。図6において、関係表現としては、言語的な意味関係を表わす接続詞や接続助詞（たとえば、「というのは」、「だって」、「（だ）から」）などが関係名毎に定義される。

【0067】ステップS209において、関係表現が記述単位Rに含まれる（すなわち、YES）と判定された場合、関連表現抽出部4は、この記述単位Rを関連表現として抽出する。

【0068】たとえば、図5において、番号「3」の記述単位R「すっきりした感じで、コクもあるからです」の中には、図6に参照される関係名「理由」に対応した関係表現「（だ）から」が存在するため、関連表現として番号「3」を抽出することになる。

【0069】次に、関連表現として抽出された記述単位Rに関係名のラベルを付与し（ステップS210）、記述単位Rを関連表現リストに追加して（ステップS211）、ステップS207に戻る。

【0070】関連表現リストは、たとえば、図7に示すような情報を持つリストとして表現される。図7において、たとえば、番号「3」の関連表現「すっきりした感じで、コクもあるからです」に対応して、関係名「理由」および位置「53-71」が対応付けられる。

【0071】上記関連表現リストの作成処理（ステップS207～S211）は、記述単位Rが取得できる限り、繰り返し実行される。その後、取得すべき記述単位Rが存在しなくなった（ステップS207において、NOと判定された）時点で、関連表現抽出部4は、関連表現リストの作成処理を終了し、特徴表現リスト（図4参照）と、関連表現リスト（図7参照）とを、関連表現統合部6に出力する。

【0072】これにより、関連表現統合部6は、関連表現統合知識記憶部5を参照して特徴表現と関連表現とを統合し（ステップS212）、関連表現統合結果として分類部7に出力する。

【0073】ここで、関連表現統合知識記憶部5は、ある特徴表現に統合し得る関連表現を判定するための関連

表現統合知識を記憶している。たとえば、以下のような判定基準を用い、各関連表現に対して統合可能性のスコアを算定する。

【0074】位置情報に関しては、関連表現の開始位置が、特徴表現の終了位置に近いほどスコアを高く算定するか、または、関連表現の終了位置が、特徴表現の開始位置に近いほどスコアを高く算定する。

【0075】同一単語に関しては、特徴表現と関連表現とで共通の単語が存在すればスコアを加算する。また、類義語に関しては、特徴表現と関連表現とで類義の単語が存在すればスコアを加算する。

【0076】単語の類義性の判定方法としては、たとえば、特徴表現および関連表現に出現する単語に関して、単語の出現頻度に基づく統計的な重み付けを行って作成した索引表を用いて、文章同士の類似度を単語ベクトルの内積として算定する方法がある。または、類義語辞書を使用する判定方法も可能である。

【0077】関連表現統合部6は、上記のように算定されたスコアの高い関連表現を特徴表現に統合し、関連表現統合結果を生成する。関連表現統合結果は、たとえば、図8に示すように、特徴表現に対して統合された関連表現を「木構造」の形式で構成したものととなる。

【0078】図8に示した関連表現統合結果の例では、抽出意図が番号「1」の「賞賛」に対し、番号「3」の関係名「理由」と、番号「4」の関係名「逆接」とが統合されている。

【0079】図2に戻り、最後に、分類部7は、関連表現統合部6から入力された関連表現統合結果を内容に応じて分類し（ステップS213）、図2の処理を終了する。このとき、分類部7は、たとえば、関連表現統合結果に付与されている抽出意図（図3、図4、図8参照）に応じて分類を行う。

【0080】上記処理手順により、分析対象に依存する概念定義やシソーラスを用いずに、表層的な表現に着目して表現抽出を行うことができ、分析対象に依存しない汎用性を有しつつ、一般のWebページも含めて分析することができる。

【0081】また、抽出対象となる表現A（クレーム、要望、賞賛などの抽出意図）に加えて、抽出した表現Aに対して関係（理由、例示、逆接などの関係名）を有する表現Bを抽出し、表現A、Bを関連付けて提示することにより、有益な情報の発見や、それらの関係の分析を支援することができる。

【0082】実施の形態2. なお、上記実施の形態1では、関連表現統合知識記憶部5内の知識のみを関連表現統合部6と関連させたが、さらに対象名辞書内の情報を関連させてもよい。

【0083】図9はこの発明の実施の形態2によるテキストマイニング装置を示すブロック構成図であり、前述（図1参照）と同様のものについては、同一符号を付し

て詳述を省略する。

【0084】8は対象名辞書であり、分析の目的とする対象名に関する情報を記憶する。9は対象名抽出部であり、対象名辞書8と関連表現統合部6との間に介在されており、対象名を含む第3の記述単位を抽出して、対象名とともに出力する。

【0085】この場合、関連表現統合部6は、統合対象となる第1の記述単位である特徴表現抽出結果（図4に参照の特徴表現リスト）と関連表現抽出結果（図7に参照の関連表現リスト）との中から、対象名が抽出された第3の記述単位に対象名を付与して分類部7に出力する。

【0086】これにより、分類部7は、関連表現統合部6から入力された関連表現統合結果に付与されている対象名毎に応じて、関連表現統合結果を分類して提示する。

【0087】また、関連表現統合部6は、特徴表現抽出結果と関連表現抽出結果との統合判定に際して、対象名を利用する。さらに、分類部7は、関連表現統合部6により得られた関連表現統合結果を事例として、事例のクラスタリングを行い、問題解決木を構成して出力する。

【0088】図9において、図1と異なる点は、対象名辞書8および対象名抽出部9を備えていることのみである。対象名辞書8に記憶される対象名としては、たとえば、企業名、製品名、機器名、部品名、人名などの固有名詞である。

【0089】対象名辞書8は、上記の固有名詞をリストとして記憶し、対象名抽出部9は、与えられた記述単位の中から対象名を抽出し、該当する記述単位の情報として対象名を付与する。

【0090】したがって、関連表現統合部6は、ある特徴表現に統合し得る関連表現を判定するための知識として、対象名抽出部9で付与された対象名を利用することにより、関連表現が含まれない記述単位Rに関しても、対象名に関する記述を有する関連情報として、特徴表現に統合することができる。

【0091】また、分類部7は、特徴表現抽出部2により抽出された特徴表現と、関連表現抽出部4により抽出された関連表現（理由や例示などの関係を有する）と、関連表現統合部6により統合された対象名とからなる組を事例として、事例のクラスタリングを行い、「問題解決木」を構成して出力することができる。なお、クラスタリング技術については、公知の方法が種々提案されており、ここでは、どの方法を用いても構わない。

【0092】上記のように、図9の構成を用いて処理することにより、分析の目的とした対象名に関する、競合製品や他社の評判、関連する機器名／部品名の関連情報を抽出することができる。

【0093】また、分類部7において、事例のクラスタリングを行うことにより、分析の目的とした対象名に関

する、不評、クレーム、障害発生などの問題に関して、共通の現象（特徴表現）に関する理由や具体例などの分析を効果的に支援することができる。

【0094】

【発明の効果】以上のように、この発明によれば、入力文書に関して、抽出意図と意図抽出表現とを対応付けるための特徴表現抽出規則を記憶する特徴表現抽出規則記憶部と、入力文書から、意図抽出表現を含む第1の記述単位を抽出して特徴表現とし、特徴表現を含む解析結果を特徴表現抽出結果として、抽出意図とともに出力する特徴表現抽出部と、第1の記述単位毎の関係に関し、関係名と関係を表わす関係表現とを対応付けるための関連表現抽出規則を記憶する関連表現抽出規則記憶部と、特徴表現抽出結果から、関係表現を含む第2の記述単位を抽出して関連表現とし、関連表現を含む解析結果を関連表現抽出結果として、関係名とともに出力する関連表現抽出部と、特徴表現と関連表現とを第1の記述単位毎の関係に基づいて統合するための関連表現統合知識を記憶する関連表現統合知識記憶部と、関連表現統合知識を参照して特徴表現と関連表現とを統合し、関連表現統合結果として出力する関連表現統合部と、関連表現統合結果を抽出意図の内容に応じて分類する分類部とを備え、概念やシソーラスを用いずに表層的な表現に着目して表現抽出を行うようにしたので、Webページも含めて分析対象にできるような汎用性を有するテキストマイニング装置が得られる効果がある。また、抽出対象となる表現A（意見、クレーム、障害報告など）に加えて、抽出した表現Aと関係（根拠や理由、原因、例示など）を有する表現Bを抽出して、表現A、Bを関連付けて提示するようにしたので、有益な情報の発見や、それらの関係の分析を支援することのできるテキストマイニング装置が得られる効果がある。

【0095】また、この発明によれば、分析の目的とする対象名に関する情報を記憶する対象名辞書と、対象名を含む第3の記述単位を抽出し、対象名とともに出力する対象名抽出部とを備え、関連表現統合部は、統合対象となる第1の記述単位である特徴表現抽出結果と関連表現抽出結果との中から、対象名が抽出された第3の記述単位に対象名を付与して出力し、分類部は、対象名毎に分類して提示するようにしたので、対象名に関する、競合製品や他社の評判、関連する機器名／部品名の関連情報を抽出することのできるテキストマイニング装置が得られる効果がある。

【0096】また、この発明によれば、分類部は、関連表現統合部により得られた関連表現統合結果を事例として、事例のクラスタリングを行い、問題解決木を構成して出力するようにしたので、対象名に関する、不評、クレーム、障害発生などの問題に関して、共通の現象（特徴表現）に関する理由や具体例などの分析を効果的に支援することのできるテキストマイニング装置が得られる

効果がある。

【0097】また、この発明によれば、入力文書に関して、抽出意図と意図抽出表現とを対応付けるための特徴表現抽出規則を記憶する特徴表現抽出規則記憶ステップと、入力文書から、意図抽出表現を含む第1の記述単位を抽出して特徴表現とし、特徴表現を含む解析結果を特徴表現抽出結果として、抽出意図とともに出力する特徴表現抽出ステップと、第1の記述単位毎の関係に関し、関係名と関係を表わす関係表現とを対応付けるための関連表現抽出規則を記憶する関連表現抽出規則記憶ステップと、特徴表現抽出結果から、関係表現を含む第2の記述単位を抽出して関連表現とし、関連表現を含む解析結果を関連表現抽出結果として、関係名とともに出力する関連表現抽出ステップと、特徴表現と関連表現とを第1の記述単位毎の関係に基づいて統合するための関連表現統合知識を記憶する関連表現統合知識記憶ステップと、関連表現統合知識を参照して特徴表現と関連表現とを統合し、関連表現統合結果として出力する関連表現統合ステップと、関連表現統合結果を抽出意図の内容に応じて分類する分類ステップとを備え、概念やシソーラスを用いずに表層的な表現に着目して表現抽出を行うようにしたので、Webページも含めて分析対象にできるような汎用性を有するテキストマイニング方法が得られる効果がある。

【0098】また、この発明によれば、分析の目的とする対象名に関する情報を記憶する対象名辞書から、対象名を含む第3の記述単位を抽出し、対象名とともに出力する対象名抽出ステップを備え、関連表現統合ステップは、統合対象となる第1の記述単位である特徴表現抽出結果と関連表現抽出結果との中から、対象名が抽出された第3の記述単位に対象名を付与して出力し、分類ステップは、対象名毎に分類して提示するようにしたので、対象名に関する、競合製品や他社の評判、関連する機器名／部品名の関連情報を抽出することのできるテキストマイニング方法が得られる効果がある。

＊

＊【0099】また、この発明によれば、分類ステップは、関連表現統合ステップにより得られた関連表現統合結果を事例として、事例のクラスタリングを行い、問題解決木を構成して出力するようにしたので、対象名に関する、不評、クレーム、障害発生などの問題に関して、共通の現象（特徴表現）に関する理由や具体例などの分析を効果的に支援することのできるテキストマイニング方法が得られる効果がある。

【図面の簡単な説明】

10 【図1】 この発明の実施の形態1によるテキストマイニング装置を示すブロック構成図である。

【図2】 この発明の実施の形態1による処理動作を示すフローチャートである。

【図3】 この発明の実施の形態1による特徴表現抽出規則の例を示す説明図である。

【図4】 この発明の実施の形態1による特徴表現リストの例を示す説明図である。

【図5】 この発明の実施の形態1による記述単位群の例を示す説明図である。

20 【図6】 この発明の実施の形態1による関連表現抽出規則の例を示す説明図である。

【図7】 この発明の実施の形態1による関連表現リストの例を示す説明図である。

【図8】 この発明の実施の形態1による関連表現統合結果の例を示す説明図である。

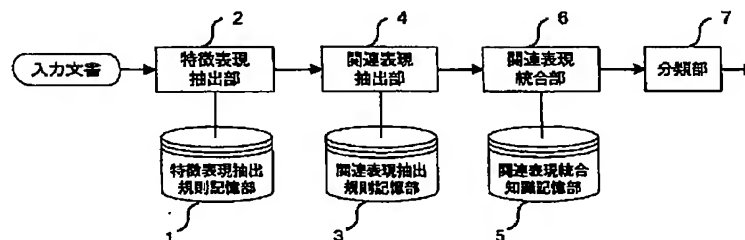
【図9】 この発明の実施の形態2によるテキストマイニング装置を示すブロック構成図である。

【図10】 従来のテキストマイニング装置を示すブロック構成図である。

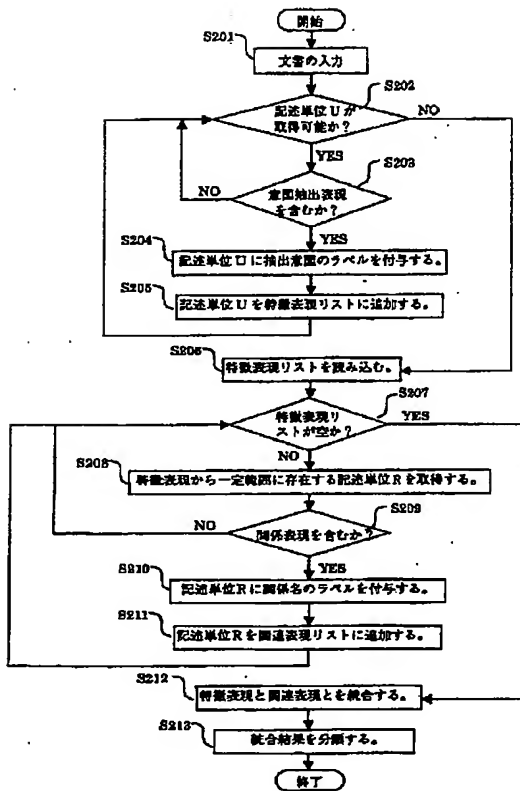
【符号の説明】

1 特徴表現抽出規則記憶部、2 特徴表現抽出部、3 関連表現抽出規則記憶部、4 関連表現抽出部、5 関連表現統合知識記憶部、6 関連表現統合部、7 分類部、8 対象名辞書、9 対象名抽出部。

【図1】



【図2】



【図3】

特徴表現抽出規則の例

抽出意図	意図抽出表現	重み
クレーム	・対応 / 悪	1.0
	・納得 / でき / ない	1.0
要望	・し / て / 欲	0.8
賞賛	・とても / 満足	1.0

【図4】

特徴表現リストの例

特徴表現	抽出意図	位置(開始-終了)	重み
「A社のサポート受付は対応が悪い」	クレーム	21 - 36	1.0
「(製品名)はすぐ壊れるのに、値段が高いのに納得できない」	クレーム	119 - 148	1.0
「講習会を開催して欲しいなあと思います」	要望	99 - 117	0.8
「この秋に出た(製品名)の味にとても満足しています」	賞賛	1 - 26	1.0

【図6】

関連表現抽出規則の例

関係名	関係表現
理由	「というのは」「だって」「(だ)から」「のため」「起因」
例示	「例えば」「一例」
順接	「それで」「そこで」「つまり」
逆接	「しかし」「でも」「ところが」「(だ)けど」
並列	「また」「さらに」「一方」
...	...

【図5】

記述単位群の例

番号	記述単位群	位置(開始-終了)
1	「この秋に出た(製品名)の味にとっても満足しています」	1 - 26
2	「これは今までの製品とはちょっと違うような気がします」	27 - 52
3	「すっきりした感じでコクもあるからです」	53 - 71
4	「でも、デザインがいまいちですね」	72 - 87
5	「まあ、いいか」	88 - 94

【図7】

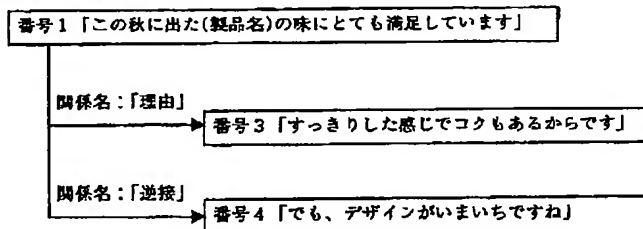
関連表現リストの例

番号	関連表現	関係名	位置(開始-終了)
3	「すっきりした感じで、コクもあるからです」	理由	53 - 71
4	「でも、デザインがいまいちですね」	逆接	72 - 87
...

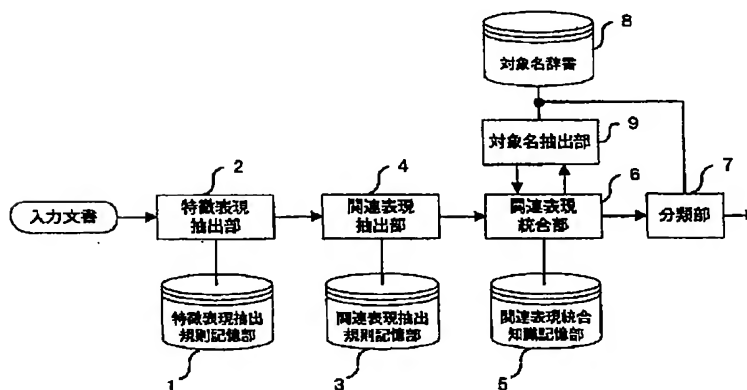
【図8】

関連表現統合結果の例

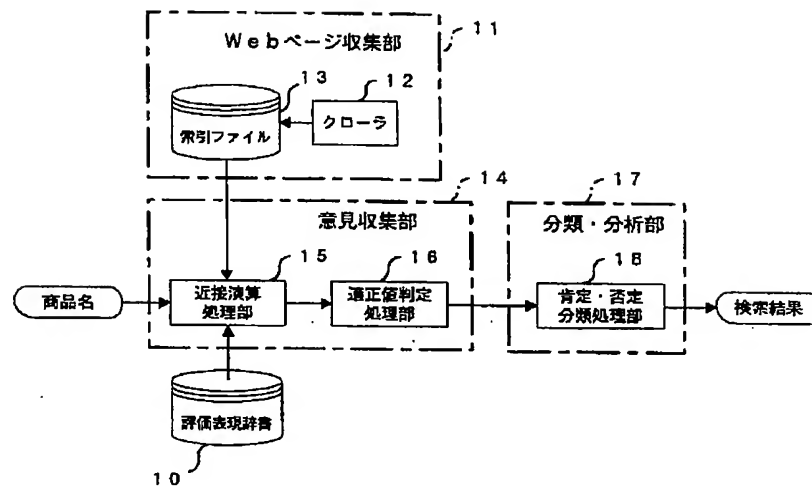
抽出意図：賞賛



【図9】



【図10】



フロントページの続き

(72)発明者 鈴木 克志
東京都千代田区丸の内二丁目2番3号 三
菱電機株式会社内

F ターム(参考) SB075 ND03 NK06 PP24 UU06